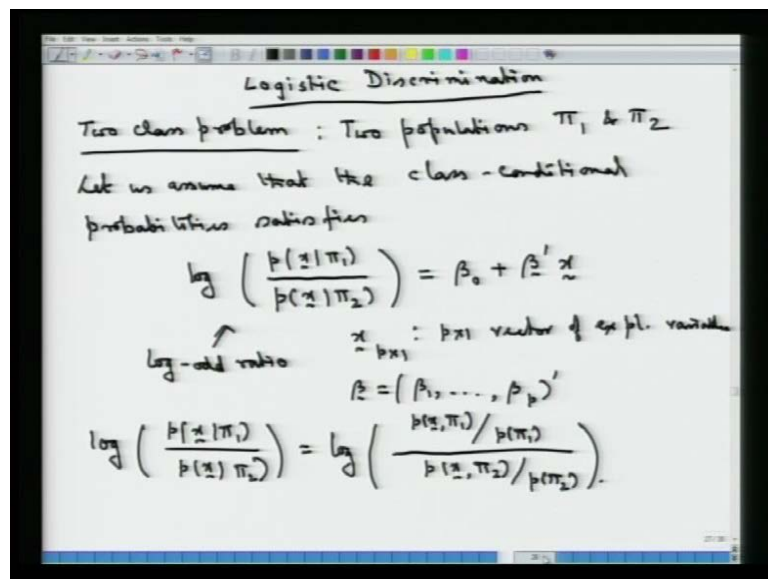


Applied Multivariate Analysis
Prof. Amit Mitra
Prof. Shramishtha Mitra
Department of Mathematics and Statistics
Indian Institute of Technology, Kanpur
Lecture No. #35
Discriminant Analysis and Classification

In this lecture, we are going to talk about logistic discrimination. So, we are going to build discrimination function that is going to be based on the principal of logistic regression. So, let us try to look at what we are up to.

(Refer Slide Time: 00:27)



So, we are talking about construction of logistic discrimination. To start with a simple problem, let us consider a two class problem meaning thereby we have got two populations; namely π_1 , and π_2 . And we assume that the class conditional probabilities satisfy the following relationship.

Let us assume that the class conditional probabilities (No audio from 01:12 to 01:21) satisfies the following relationship that, \log of $p(x|\pi_1)$. So, this is the density, this is conditioned on the fact that x is coming from π_1 , that divided by the density of x ,

Now we had earlier denoted this $p_{i|j}$ to be equal to p_i simply. So, for notational convenience, we are writing this as p_i of $p_{i|j}$ equal to p_i . So, this log odds ratio, thus is going to be equal to these two terms cancel out, and what we will be having the log odds ratio as this is a p_1 times $p_{i|1}$ given x that divided by p_2 into $p_{i|2}$ given x .

And from the condition of the logistic discrimination, this is equal to our beta naught plus beta prime times x . Now, this further would imply that this ratio of this posterior probabilities $p_{i|1}$ given x , this divided by $p_{i|2}$ given x quantity that is going to be equal to p_1 by p_2 times e to the power beta naught plus a beta prime x . Let us write this in the form that it is log of p_1 by p_2 term, this plus this beta naught plus a beta prime x quantity. And let us write this further as e to the power a constant beta star which is beta star is beta naught plus log of p_1 by p_2 , this plus beta prime x wherein we have used the fact that is beta star, we are denoting this as beta naught plus log of p_1 by p_2 .

So, this implies what we have here, now note that $p_{i|1}$ given x plus $p_{i|2}$ given x that is equal to 1. So, what we can write in place of $p_{i|1}$ given x is 1 minus $p_{i|2}$ given x that divided by $p_{i|2}$ given x . So, that is equal to e to the power beta star plus a beta prime times x . Now, either you write it in this particular form or one can write this as e to the power a vector beta star multiplied beta star prime that multiplied by x star, wherein what would be this beta star vector? The beta star vector would have a beta star in the first place and then that augmented by this old beta vector, and x star similarly is going to be equal to one augmented with this x vector, the original feature vector.

Because this one is going to get from this beta star and lead us to this beta star plus beta prime x . Now from this expression, what we have in here. It is further that what we have is one can express this ratio in terms of just $p_{i|2}$ given x .

(Refer Slide Time: 07:12)

Handwritten mathematical derivation on a whiteboard:

$$\Rightarrow P(\pi_2 | x) = \frac{1}{1 + e^{\beta_0^* + \beta_1' x}} \left(= \frac{1}{1 + e^{\beta_0^* + \beta_1' x}} \right)$$

$$\Delta P(\pi_1 | x) = 1 - P(\pi_2 | x) = \frac{e^{\beta_0^* + \beta_1' x}}{1 + e^{\beta_0^* + \beta_1' x}}$$

Assignment rule:

Assign x to π_1 if $\frac{P(\pi_1 | x)}{P(\pi_2 | x)} > 1$

Δ x to π_2 if \forall .

i.e. assign x to π_1 if $\beta_0^* + \beta_1' x > 0$

Δ x to π_2 if \forall .

So, this is going to imply that $p(\pi_2 | x)$ is going to be just equal to $1 / (1 + e^{\beta_0^* + \beta_1' x})$. Let us still write this as a $\beta_0^* + \beta_1' x$ quantity, because in other two just have some understanding that this is what is corresponding to the constant term out here.

So, let us just keep this as a $\beta_0^* + \beta_1' x$ quantity. So, that we have this here. So, its $\beta_0^* + \beta_1' x$ this plus $\beta_1' x$ quantity either in this form or in the form of $1 + e^{\beta_0^* + \beta_1' x}$ to the power $\beta_0^* + \beta_1' x$, which are equivalent. So, when we have this posterior probability of the second population given x to be given by this, one can also see that what $p(\pi_1 | x)$ is. So, that is going to be $1 - p(\pi_2 | x)$, and that from this expression is just going to be equal to $e^{\beta_0^* + \beta_1' x} / (1 + e^{\beta_0^* + \beta_1' x})$.

So, these are the two posterior probabilities of the respected population. So, we have obtained that it is $p(\pi_1 | x)$ is given by this and $p(\pi_2 | x)$ is given by this. So, the assignment rule or the **discriminate** discrimination rule is the following assign x to π_1 . If the posterior probability of π_1 population is higher than that of the second population, if that is assign x to π_1 if we have got $p(\pi_1 | x)$ greater than $p(\pi_2 | x)$ given x or in terms of this odds ratio what we can say is that if this ratio is greater than 1 and x to π_2 if otherwise.

So, that is basically the assignment rule that we get if we look at such a formulation that the log odds ratio satisfies that particular relationship between the log odds and that of the explanatory variables, the feature vector which is contained in x . Now a thing that should be noted here at this particular point of time, that when we are trying to implement this particular assignment rule in practice, this $p(\pi_1 | x)$ or $p(\pi_2 | x)$ depends on the parameters β_0 . And the parameters which represent in this beta vector which is $\beta_1, \beta_2, \dots, \beta_p$.

And from the given data, from the learning sample one has to thus at that particular stage come up with some estimates of these quantities. And based on the estimated these quantities will be actually implementing this classification rule. Now let us try to extend this approach, now it is another way to look at this particular term is one can see that, that is assign x to π_1 if **if** we are looking at what is this ratio $p(\pi_1 | x)$ by $p(\pi_2 | x)$, it is this divided by this particular quantity. So, it is going to be $e^{\beta_0 + \beta_1'x}$ that is greater than 1. Or in other words we will have this rule to be given by this $\beta_0 + \beta_1'x$, that to be greater than $\log(1/0)$.

And x to π_2 , if we have the condition otherwise, that is in a more simple manner it is this straight forward depending on these unknown quantities β_0 plus 1 of them a β_1 and a β_2 quantities which are there in this beta vector.

(Refer Slide Time: 11:48)

Multi-class problem

C pop^{ns} : π_1, \dots, π_C .

Assume the the log-odds for any pair satisfies.

$$\log \left(\frac{p(\pi_s | x)}{p(\pi_c | x)} \right) = \beta_{s0} + \beta_s' x \quad ; s=1, \dots, C-1$$

i.e. $C-1$ log-odds specify the model
using these $C-1$ log-odds.

$$p(\pi_s | x) = \frac{\exp(\beta_{s0} + \beta_s' x)}{1 + \sum_{s=1}^{C-1} \exp(\beta_{s0} + \beta_s' x)}$$

$s=1, \dots, C-1$

So, we now look forward to extending this particular logistic discrimination function in a c class problem a multiclass problem, wherein we assume that we have got c populations multiclass problem. And trying to frame, what is going to be the logistic discrimination rule in such a situation.

So, we assume that there are c populations, which are say denoted by $\pi_1 \pi_2 \dots \pi_c$. Now what we are now going to assume is the following, assume that the log odds for every pair satisfies any pair of populations, satisfies the following relationship that the log of $p_{\pi_i} | x$ given π_i say s th population, that divided by $p_{\pi_c} | x$. This satisfies an equation say is equal to β_s naught plus a β_s vector transpose x , this is for s equal to 1 to up to c . So, we are looking at, this is the pair of π_c th population and any of the π_i s population for i equal to 1 up to. I am **sorry** and for s equal to 1 to up to c minus 1. So, we are looking at pairs of such populations, and then looking at the log odds of these quantities which we assume that it is of this particular form. So, this would imply that this c minus 1 log odds specify the model null. So, we have got this c minus 1 log odds, which we get from this expression for s equal to 1 to up to c minus 1, these c minus 1 log odds specify the model completely. Now, using these c minus 1 log odds and using the fact that summation of all these quantities $p_{\pi_i} | x$ summation from i equal to 1 to up to c , these are going to add up to one. What we can now see is the following, using this c minus 1 log odds, what will be getting is that $p_{\pi_i} | x$ the posterior probabilities $p_{\pi_i} | x$ given x , this is going to be given by e to the power β_s naught star, it is almost the same as what we had for the two class problem. That plus β_s transpose x that divided by 1 plus this summation of s equal to 1 to up to c minus 1 of these quantities e to the power β_s naught star, this plus β_s prime times x , and this is for s equal to 1 to up to c minus 1.

(Refer Slide Time: 15:26)

$$p(\pi_c | x) = \frac{1}{1 + \sum_{s=1}^{c-1} e^{x(\beta_{s0}^* + \beta_s' x)}}$$

where, $\beta_{s0}^* = \beta_{s0} + \log\left(\frac{p(\pi_s)}{p(\pi_c)}\right)$

Assignment rule:
 Assign x to π_k if $p(\pi_k | x) = \max_i p(\pi_i | x)$

So, for these $c - 1$ populations, the posterior density $p(\pi_s | x)$ is going to be given by this. I will just define, what is β_{s0}^* and those quantities. And you will have this $p(\pi_c | x)$, that to be given by 1 upon the same denominator. So, this is 1 upon 1 plus summation s equal to 1 to up to $c - 1$, this also was $c - 1$ of e to the power the same quantity as before. So, that was denoted by β_{s0}^* , this plus this $\beta_s' x$ quantities. So, it is this term, wherein we have got as in the previous setup this β_{s0}^* , that is going to be given by β_{s0} the first one, then that multiplied by \log of $p(\pi_s)$ which we have denoted by p_s quantity simply that divided by $p(\pi_c)$.

Now, once we have these expressions, c of these out here. This is the posterior density of the s th population given x , this is for s equal to 1 to up to $c - 1$, and this is for the c th population. Assignment rule is the following, using this logistic discrimination function is just the extension of what we have for the two class problem. Assign x to π_k if we have $p(\pi_k | x)$ the posterior probability of the k th population given x , if that is maximum over i of all these posterior probabilities $p(\pi_i | x)$, given x right. So, this is what we have as the classification rule, we are going to assign x to π_k , if we have got this to be true.

Now, this once again can be expressed in terms of this β_{s0}^* and β_s vectors, because all these are that.

(Refer Slide Time: 17:51)

Class-conditional probabilities & expected

Consider Y to be binary, 1 and 0 (2 class prob).

$$Y = \begin{cases} 1 & \text{if } \pi_1 \\ 0 & \text{if } \pi_2 \end{cases}$$

We have seen that

$$P(Y=1|x) = P(\pi_1|x) = \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)}$$
$$P(Y=0|x) = P(\pi_2|x) = \frac{1}{1 + \exp(\beta^T x)}$$
$$\Rightarrow E(Y|x) = \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)} = P(\pi_1|x)$$

So, one can find out what is that quite easily. Now let us look at the following relationship, which is interesting the name why logistic discrimination comes at all, let us now look at the class conditional probabilities **class conditional probability** and expected response. So, what is the relationship between these, and how we are going to get the name logistic discrimination? Now for simplicity, consider a two class problem consider a Y a variable to be binary say taking value 0 and 1. So, this is corresponding to the two class problem.

Now, in what sense is this two class problem. So, Y takes the value say 1, if the population is π_1 and takes the value 0, if the population is π_2 . So, corresponding to the particular feature vector, we of course have the in the pre classified examples in the learning sample, what is the class membership of that particular feature vector. And if that is π_1 , then the value of the Y variable which is a binary variable. We take that to be equal to 1, and if the membership of the feature vector is π_2 population, then we take the value of this binary variable to be equal to 0.

Now, what we have seen earlier is that, we have already seen that, now these are that is two class problems. Now we are going to write this expression of $p(\pi_1|x)$ which we have discussed today, which is e to the power. Let us look back and see what it is for the two class problem, we had this $p(\pi_1|x)$ given x to be given by this expression, and $p(\pi_2|x)$ given x to be given by this expression. So, let us use those expressions and

write this as a $\beta'x$ divided by $1 + e^{\beta'x}$.

And this quantity, if you now look at this binary random variable Y . So, this has got two values 1 and 0 corresponding to this. So, if we are looking at p_1 given x , then that is nothing, but in terms of the Y variable, its probability that Y is equal to 1 given x . So, this is what is a conditional probability of Y taking the value 1 given x is observed, and similarly we have also seen that p_2 given x that is equal to 1 upon the same denominator. So, that is $1 + e^{\beta'x}$.

Now, what is this equal to this in terms of the binary variable. This is probability that Y is equal to 0 given x . So, we if we have got these two as the conditional probability masses of Y taking the value 1, and Y taking the value 0. This would imply that if we look at the conditional expectation of Y given x , what is that going to be equal to? This takes the value 1 with this probability and 0 with this probability. And hence, the conditional expectation of Y given x is just going to be given by this particular expression, which is $e^{\beta'x}$ divided by $1 + e^{\beta'x}$, which by the way is nothing, but p_1 given x .

Now, from this expression here, which we have got this conditional expectation to be equal to this term here. Let us now denote this quantity $\beta'x$ to be equal to a quantity which is θ , denote by θ the quantity which is $\beta'x$.

(Refer Slide Time: 22:04)

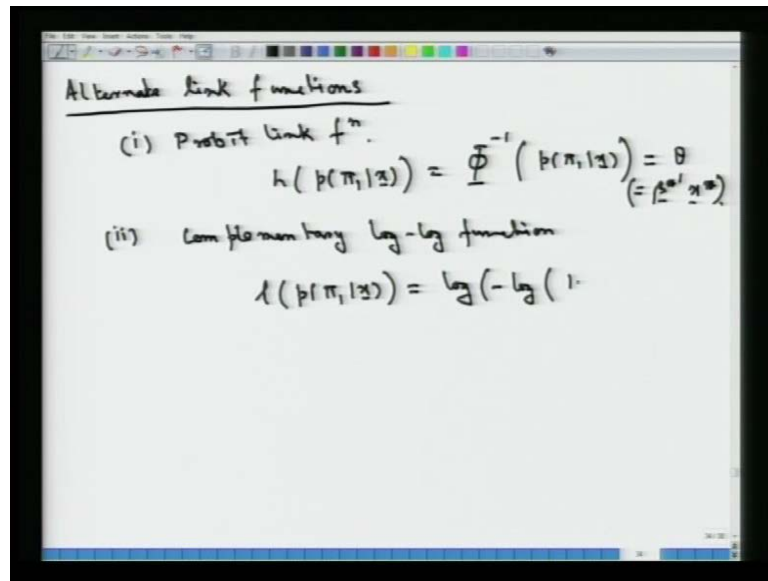
$$\begin{aligned} \text{Denote by } \theta &= \beta' x^* \\ \Rightarrow p(\pi_1 | x) &= \frac{e^\theta}{1 + e^\theta} = \frac{1}{1 + e^{-\theta}} \\ \Rightarrow 1 + e^{-\theta} &= \frac{1}{p(\pi_1 | x)} \\ \Rightarrow e^{-\theta} &= \frac{1}{p(\pi_1 | x)} - 1 \\ \text{i.e. } \theta &= \log\left(\frac{p(\pi_1 | x)}{1 - p(\pi_1 | x)}\right) = g(p(\pi_1 | x)) \\ &\quad \uparrow \\ &\quad \text{logistic link f}^\bullet \end{aligned}$$

So, this would imply that what we have as $p(\pi_1 | x)$ given x is nothing, but equal to, let us see this expression it is going to be equal to e to the power θ divided by $1 + e$ to the power θ . So, that is equal to e to the power θ that divided by $1 + e$ to the power θ , which one can also write as one upon $1 + e$ to the power minus θ .

So, if we have this expression, this would imply that $1 + e$ to the power minus θ that is equal to one upon $p(\pi_1 | x)$ given x . So, this would imply that e to the power minus θ is equal to one upon $p(\pi_1 | x)$ given x , this minus 1. That is, one can write the θ quantity in terms of log of this other way round. So, it is going to be equal to log of $p(\pi_1 | x)$ given x this divided by $1 - p(\pi_1 | x)$ given x , now this is some function of $p(\pi_1 | x)$ given x .

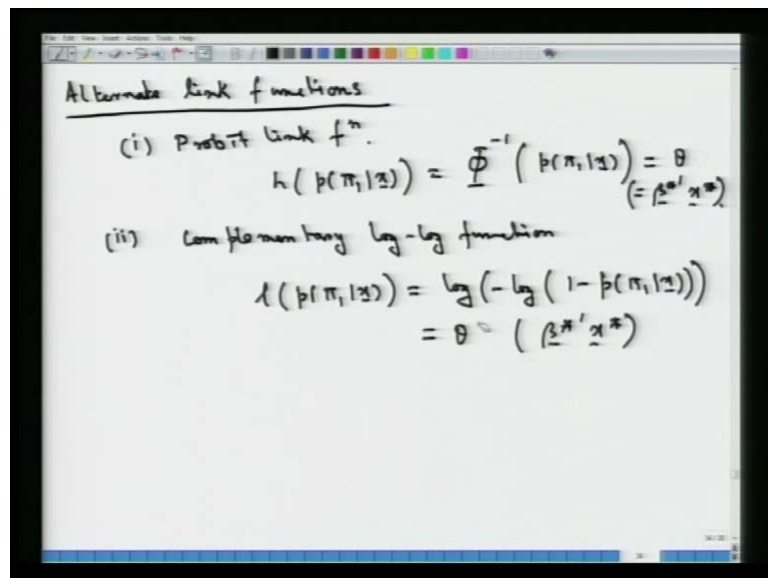
Now, since we have got this function form in terms of this log odds ratio here, and that $p(\pi_1 | x)$ given x is of this logistic function. We have the name that, this link that we are linking θ with $p(\pi_1 | x)$ given x the posterior density of π_1 given x through a logistic link function. So, this is the logistic link function, and hence what we get the name as the logistic discrimination. Now alternate forms of having this link function, because this is just the relationship between this θ , which is $\beta' x^*$ and this $p(\pi_1 | x)$ given x .

(Refer Slide Time: 24:36)



So, if this function is of the form that it is log of p pi 1 given x or rather it is assumed to be of the form that it is log of p pi 1 given x, that divided by 1 minus p pi 1 given x. Then what we get is the logistic discrimination alternate link functions are following.

(Refer Slide Time: 24:35)



Alternate link functions, there are two popular alternate link functions. One which is called the profit link, profit link function wherein we assume that this function, what we have a function of p pi 1 given x this is a probability. So, this h of this p pi 1 x is going to be given by capital phi inverse of p pi 1 given x.

So, this particular function here phi inverse, where capital phi is the probability distribution function of a standard normal variate that is equal to theta, which is our $x^T \beta$. So, this is equal to our $x^T \beta$. Now the second type of popular link function is what is called the complimentary log function, a complimentary log function **log** function rather, wherein we assume that we have got a function of this $p(\pi_1 | x)$ given x to be given by the following quantity which is \log of $\frac{p(\pi_1 | x)}{1 - p(\pi_1 | x)}$.

So, that this bracket ends here and that is linked with this theta, which once again is that linear combination of the parameters with the feature vector with the constant vector 1 attached to it.

(Refer Slide Time: 26:28)

Define by $\theta = \beta^T x^*$

$$\Rightarrow p(\pi_1 | x) = \frac{e^\theta}{1 + e^\theta} = \frac{1}{1 + e^{-\theta}}$$

$$\Rightarrow 1 + e^{-\theta} = \frac{1}{p(\pi_1 | x)}$$

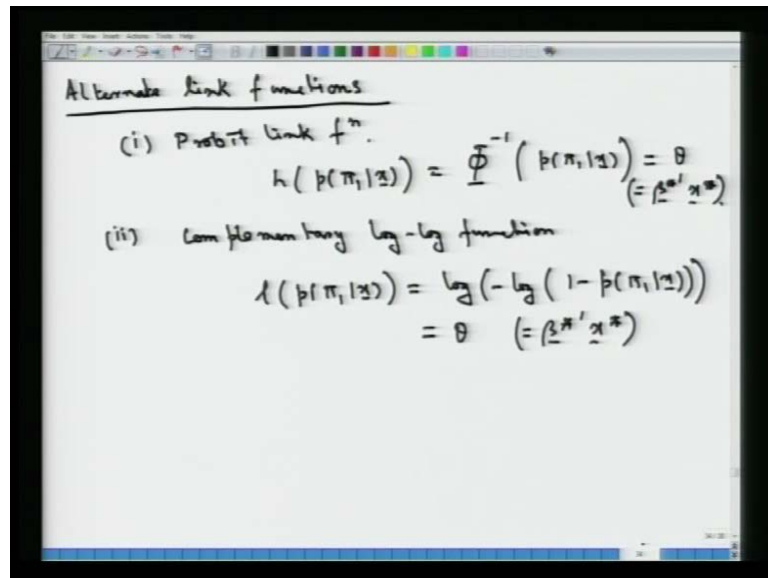
$$\Rightarrow e^{-\theta} = \frac{1}{p(\pi_1 | x)} - 1$$

$$\text{i.e. } \theta = \log \left(\frac{p(\pi_1 | x)}{1 - p(\pi_1 | x)} \right) = g(p(\pi_1 | x))$$

↑
Logistic link f^*

So, that these are two alternate link functions. In case of a logistic discrimination, the link that we have is precisely this that the functional link between $p(\pi_1 | x)$ and theta is given by this expression which is logistic link, this is the profit link. And this is the complimentary log **log** function link.

(Refer Slide Time: 26:40)



Alternate link functions

(i) Probit link f^n .

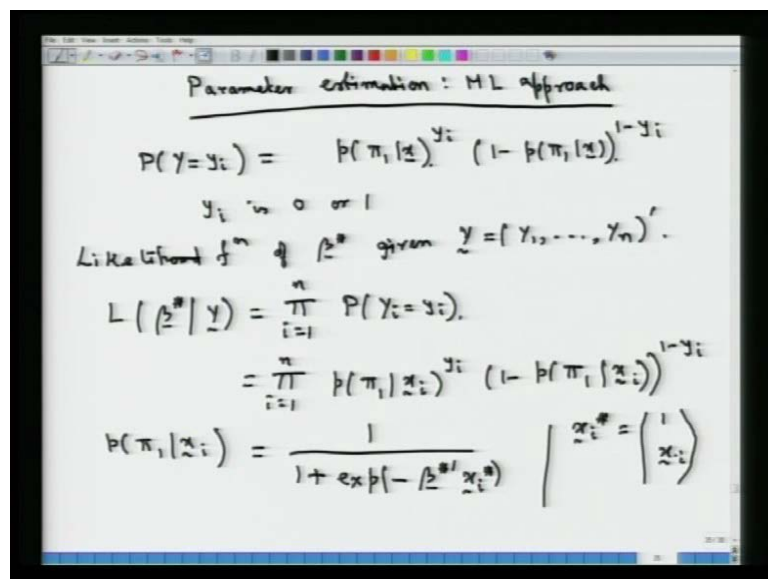
$$h(p(\pi_i|z)) = \Phi^{-1}(p(\pi_i|z)) = \theta \quad (= \beta^{*T} z^*)$$

(ii) Complementary log-log function

$$l(p(\pi_i|z)) = \log(-\log(1-p(\pi_i|z))) = \theta \quad (= \beta^{*T} z^*)$$

Now let us last look at or rather talk about little bit about parameter estimation.

(Refer Slide Time: 26:52)



Parameter estimation: ML approach

$$P(Y=y_i) = p(\pi_i|z)^{y_i} (1-p(\pi_i|z))^{1-y_i}$$

y_i is 0 or 1

Likelihood of β^* given $y = (y_1, \dots, y_n)^T$.

$$L(\beta^*|y) = \prod_{i=1}^n P(Y_i=y_i)$$
$$= \prod_{i=1}^n p(\pi_i|z_i)^{y_i} (1-p(\pi_i|z_i))^{1-y_i}$$
$$p(\pi_i|z_i) = \frac{1}{1 + \exp(-\beta^{*T} z_i^*)} \quad \left| \quad z_i^* = \begin{pmatrix} 1 \\ z_i \end{pmatrix} \right.$$

Now the method of parameter estimation that is usually adapted in such situation is the method of maximum likelihood. So, parameter estimation M L approach let us still look at that simple formulation that we have got two possibilities 0 and 1. So, we have got under such a situation of the previous setup, what we have discussed here in this formulation that we have got Y_i es to be defined in this particular way, and in such a situation what we are going to have is the following.

Probability that Y is equal to y_i that is going to be given by $p^{\pi_i} (1-p)^{1-\pi_i}$ given x , this to the power y_i into $1 - p^{\pi_i} (1-p)^{1-\pi_i}$ given x , this to the power $1 - y_i$. Now, y_i takes either of the two values y_i is 0 or 1. That is we are looking at this of course, is the conditional quantity conditional probability mass function given x . So, this is y equal to 0 is what we will have as $1 - p^{\pi_i} (1-p)^{1-\pi_i}$ given x that is $p^{\pi_i} (1-p)^{1-\pi_i}$ given x . And probability that y is equal to 1 given x that is equal to $p^{\pi_i} (1-p)^{1-\pi_i}$ given x . So, this is that particular quantity.

Now, hence if we have y_1, y_2, \dots, y_n the likelihood function, the likelihood function of the parameters, which basically is coming in that beta star vector, given this y vector equal to y_1, y_2, \dots, y_n , this is for n random samples. We are going to have this as say we denote this by $l(\beta^* \text{ vector}, \text{ this given } y \text{ vector})$. This is going to be equal to, because all the y is that is what we have they are independent. So, this is going to be the product of i equal to 1 to up to n , probability that y_i is equal to small y_i .

So, what y_i is basically denoting the i th record in the data, and that in the i th record in the data, we will have the corresponding feature vector to be denoted by x_i . So, this is going to look like the following, that it is product of i equal to 1 to up to n the product of these quantities keeping in mind that, when we are looking at probability that y_i is equal capital y_i is equal to small y_i , this is $p^{\pi_i} (1-p)^{1-\pi_i}$ given x_i vector that to the power y_i into $1 - p^{\pi_i} (1-p)^{1-\pi_i}$ given x_i that to the power $1 - y_i$. **am sorry** this is not $p^{\pi_i} (1-p)^{1-\pi_i}$ given x_i that to the power $1 - y_i$.

Now, we use the fact that this **this** $p^{\pi_i} (1-p)^{1-\pi_i}$ given x_i is nothing, but this is going to be of the form that it is in the form $1 + e^{-\beta^* \text{ transpose into } x_i \text{ star}}$ quantity, wherein $x_i \text{ star}$ is the following term $x_i \text{ star vector}$, is the vector which is one in the first entry. And then we have x_i the feature vector to make up the rest of the p entries in here. So, using this particular fact, we have this and accordingly if we plug-in the values of $p^{\pi_i} (1-p)^{1-\pi_i}$ given x_i in this expression here. We have the explicit form of this likelihood.

(Refer Slide Time: 31:06)

The image shows a handwritten derivation on a whiteboard. At the top, the likelihood function is given as:
$$\Rightarrow L(\beta^* | y) = \prod_{i=1}^n \left(\frac{1}{1 + \exp(\beta^{*T} x_i)} \right)^{y_i} \left(1 - \frac{1}{1 + \exp(\beta^{*T} x_i)} \right)^{1 - y_i}$$
Below this, the text 'log Likelihood' is written. The log-likelihood function is then derived as:
$$\ell(\beta^*) = \log L(\beta^* | y)$$

$$= \sum_{i=1}^n y_i \log \left(\frac{p(\pi_i | x_i)}{1 - p(\pi_i | x_i)} \right)$$

$$+ \sum_{i=1}^n \log (1 - p(\pi_i | x_i))$$

So, this will imply that this $1/\beta^*$ given this y vector that is going to be equal to product i equal to 1 to up to n , and then we have $p(\pi_i | x_i)$ given x . So, that now we are writing that as one upon $1 + e$ to the power $\beta^* x_i$ quantity. So, this is $p(\pi_i | x_i)$ given x that to the power y_i that multiplied by $1 - p(\pi_i | x_i)$. So, it is $1 - p(\pi_i | x_i)$ one upon $1 + e$ to the power the same quantity. So, just erase this 1 here. So, that we have a big denominator coming up. So, its $\beta^* x_i$ vector that to the power $1 - y_i$.

Now, this if this is the likelihood, one can also write the log likelihood. The log likelihood function, let us denote that by small ℓ β^* vector, which is log of this $1/\beta^*$ given y expression. And that can be written in terms of this compact summation, which I will just write in here which is summation i equal to 1 to n y_i times log of in term I am just keeping it in terms of $p(\pi_i | x_i)$. So, that the expression is not too messy, this is going to be given by the following quantity which is this log of $p(\pi_i | x_i)$ given x_i , this is combining the second term also. So, this term plus summation i equal to 1 to up to n of log of $1 - p(\pi_i | x_i)$ term.

(Refer Slide Time: 33:24)

Handwritten notes on a whiteboard:

Likelihood eqⁿs

$$\frac{\partial \log L}{\partial \beta^*} = \sum_{i=1}^n x_i^{*T} (y_i - p(\pi_1 | x_i))$$

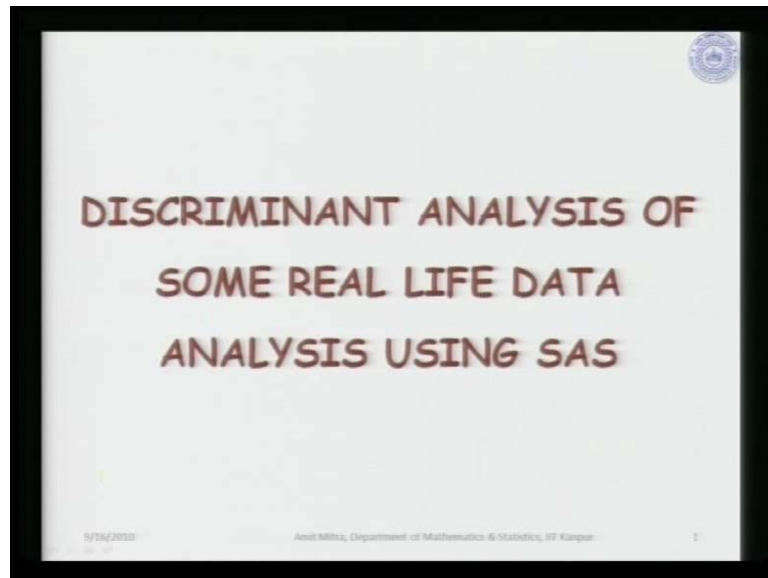
↓ System of (p+1) non-linear equations.

A method of Iteratively Reweighted Least Squares (IRLS) is applied in order to get ML estimates.

So, this is the log likelihood, then we can get to the likelihood equations from here likelihood equations corresponding to the p plus 1 parameters. We can write that compactly in the following form that this del log l with respect to this beta star vector is going to be equal to x i star transpose that into y i minus p pi 1 given x i expression. Now, this is going to be a system of non-linear equations. So, this is what we are going this summation is over i equal to one to up to n. So, this is a system of p plus 1 non-linear equations, and of course no closed form solution exists for such a system of non-linear equations. A method which is called an iteratively reweighted least squares is applied, a method of iteratively reweighted least squares or I R L S is usually applied in order to get the maximum likelihood estimates.

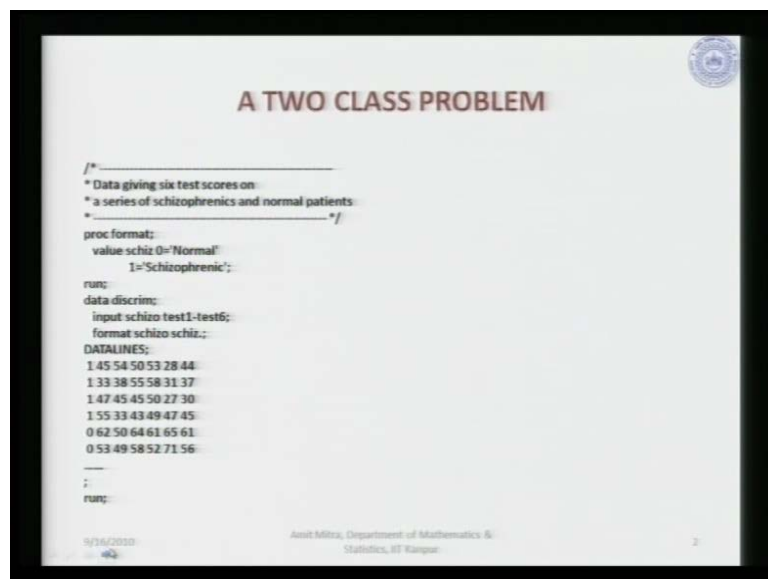
So, technically using this system of p plus 1 non-linear equations, and using I R L S, one gets to the maximum likelihood estimate, estimates of the p plus one unknown quantities. Now, once you have p plus 1, unknown quantities estimated, then one can actually look at implementation of this entire logistic discriminate function. Now, what we are next going to do is we are going to look at some real life data, and we are going to apply the type of discrimination analysis methods that we have learnt in the theoretical classes, in order to see what sort of discrimination we get in practice.

(Refer Slide Time: 35:50)



So, we now look at this a small presentation which is on a power point.

(Refer Slide Time: 36:07)



So, we are going to now look at, some real life data and real life data analysis wherein we are going to implement the type of methods that we have learnt in discriminate analysis. Now, the data analysis is done using a SAS routine. Let us look at the first example, in the first example we have a two class problem. So, there are two populations, now the two populations are following that it is a set of patients. The first

populations of patients are normal patients and the other type of population is the type of patients which are schizophrenic.

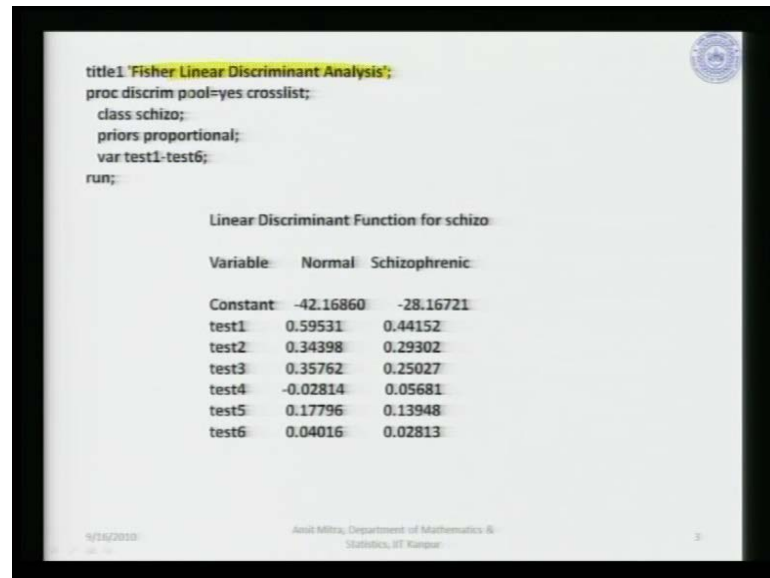
So, we have got these two classes, now we have the data, the data is giving six test scores on the series of schizophrenic and normal patients. Now, this is the format of the data, it is a huge data. So, dots are given. So, it is basically this type of data. So, it represents what the data represents is that the first entry in each of the rows this quantity here, it is basically the class identification. So, this is the data which is **the which is** comprising of the learning set data. And it has got class identification tags and hence this row of the records is what is corresponding to a schizophrenic pair patient.

Now, we are denoting that 1 if the patient is schizophrenic, and 0 if he is normal. So, we have got this to be the class membership and next six entries this 45, 54, 50, 53, 28 and 44 are the test scores corresponding to that particular patient. So, we have records like that in the data. So, some of them has class identification 1, some of them has class identification 0. So, these are normal individuals, and these are schizophrenic individuals, and these are the corresponding test scores.

Now, this six dimensional vector of test scores is now going to correspond to what we have as a feature vector. Now given this limiting sample to us, we would construct discriminate functions and classification rules based on this discriminant functions, such that we will once again look back at the learning sample. And then see how that constructive discriminant function is able to classify the pre classified examples, **how it** how it is performing on the data that is what we have in the learning sample. And then, that would lead us to the desired classification functions.

Now, the in the first case, we apply a type of discriminant function that we had studied in the very first lecture in discriminant analysis, which is a fisher linear discriminate function.

(Refer Slide Time: 38:45)



```
title1 'Fisher Linear Discriminant Analysis';
proc discrim pool=yes crosslist;
class schizo;
priors proportional;
var test1-test6;
run;
```

Linear Discriminant Function for schizo

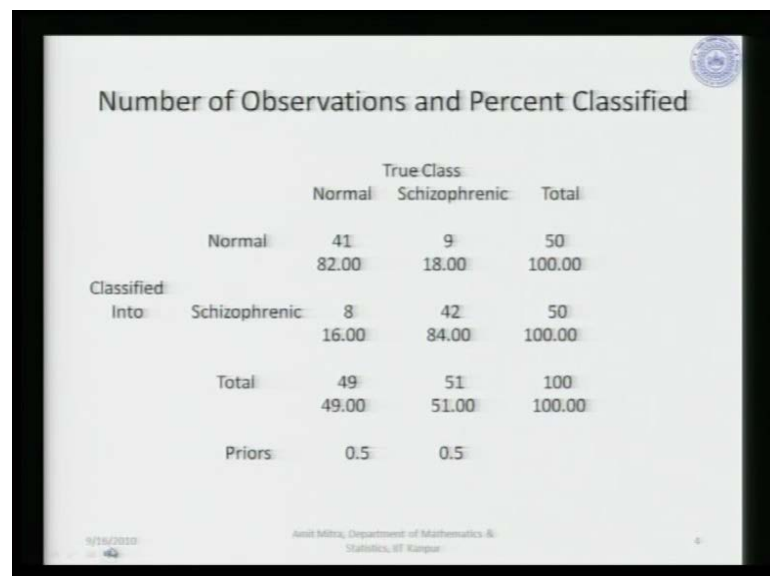
Variable	Normal	Schizophrenic
Constant	-42.16860	-28.16721
test1	0.59531	0.44152
test2	0.34398	0.29302
test3	0.35762	0.25027
test4	-0.02814	0.05681
test5	0.17796	0.13948
test6	0.04016	0.02813

9/16/2010 Amit Mitra, Department of Mathematics & Statistics, IIT Kanpur 3

So, it is implemented using a SAS code, which uses the procedure disc rim, and then looks at this particular data. Once we apply the fisher linear discriminate function to the data that we have, we get the linear discriminate function scores here, which are the coefficients that we are going to get, because this is what we have the constant term, and these are the coefficients corresponding to each of the feature vectors.

Feature vectors components are test 1, test 2, test 3, test 4, 5 and 6. And this is for the schizophrenic patients.

(Refer Slide Time: 39:18)



Number of Observations and Percent Classified

Classified Into	True Class		Total
	Normal	Schizophrenic	
Normal	41 82.00	9 18.00	50 100.00
Schizophrenic	8 16.00	42 84.00	50 100.00
Total	49 49.00	51 51.00	100 100.00
Priors	0.5	0.5	

9/16/2010 Amit Mitra, Department of Mathematics & Statistics, IIT Kanpur 4

So, these are the two different functions that we are going to get the function coefficients for the two types of patients. Now, using the fisher linear discriminate function, we look at what type of or what is the number of observations. In the learning sample, how is the performance on that learning sample of the constructed discriminant function. Now, what we get here is that now this table here is what is the confusion matrix and along with the confusion matrix, this also gives us the percent of observations which are correctly or wrongly classified.

Now, on this side we have got here. So, we have got the true class membership. So, the true class is either its normal or it is schizophrenic, and that true class observation where it is testing classified into. So, an observation coming from normal class can get classified either to the normal class or it can get classified to the schizophrenic class in which case it is going to be a misclassification. So, what we have from the given data and fisher linear discriminate function is that 41 cases, which had a true class membership of normal are now classified as normal. What is the total number of such normal patients in a class of now, total data sizes 100, among that 49 are normal patients normal individuals rather and 51 individuals are **are are** schizophrenic.

So, from among 49 normal individuals true class membership number 49, we have 41 of those been classified in to the normal category. So, this is a correct classification of the normal category individuals. Among those 49 individuals, 8 individuals among the 49 normal individuals 8 have been misclassified in to coming from the schizophrenic class. So, these are misclassifications. Now, if we look at the other class, true class membership is schizophrenic. There are 51 such patients, now from among those 51 patients; we are classifying 42 of them in to the class which is schizophrenic. And hence we are what we are doing here is a correct classification. And from among this 51 schizophrenic patients, 9 of them are classified as coming from a normal category of patients, normal category of individuals rather and hence this is misclassification.

So, this is what is giving us the confusion matrix, after we apply the fisher linear discriminant function. Herein, these two are the correct classifications, and these two are the wrong classifications.

(Refer Slide Time: 42:50)

	Normal	Schizophrenic	Total
Rate	0.1800	0.1600	0.1700
Priors	0.5000	0.5000	17%

So, we have got 18 percent of a total of 50 classifications made in the wrong category and 16 percent here in this row here going to the wrong category. And hence, now we assume here equal priors that it is 0.5, 0.5 for the two populations. Now the error counts, thus from the 2 classes pulled up its 18 percent from the normal class, 16 percent from the schizophrenic class. And hence, it is basically this percent of the observations that is 17 percent of the observations are wrongly classified using this classification rule. That is, it is simple to see that out of 100, these 17 cases are misclassified.

So, that is what is corresponding to a fisher linear discriminant function. Now corresponding to the fisher linear discriminant function, if one looks at the posterior probability of membership in to the class schizophrenic.

(Refer Slide Time: 42:52)

• Posterior Probability of Membership in schizo

Obs	From	Classified into		Normal		Schizophrenic	
		From	Into	Normal	Schizophrenic	Normal	Schizophrenic
1	Schizophrenic	Schizophrenic	*	0.1554	0.8446		
2	Schizophrenic	Schizophrenic		0.0120	0.9880		
3	Schizophrenic	Schizophrenic		0.0857	0.9143		
4	Schizophrenic	Schizophrenic		0.2352	0.7648		
5	Normal	Normal		0.9512	0.0488		
6	Normal	Normal		0.8529	0.1471		
7	Schizophrenic	Normal	*	0.9902	0.0098		
8	Schizophrenic	Schizophrenic		0.2836	0.7164		
9	Normal	Normal		0.8198	0.1802		
10	Schizophrenic	Schizophrenic		0.0022	0.9978		
11	Normal	Normal		0.8553	0.1447		
12	Normal	Normal		0.9387	0.0613		
13	Normal	Normal		0.8343	0.1657		
14	Normal	Schizophrenic	*	0.3566	0.6434		
15	Schizophrenic	Normal	*	0.8008	0.1992		
16	Normal	Normal		0.5782	0.4218		
17	Schizophrenic	Schizophrenic		0.0075	0.9925		
18	Normal	Normal		0.9498	0.0502		
19	Normal	Normal		0.9823	0.0177		
20	Schizophrenic	Normal	*	0.6446	0.3554		
21	Schizophrenic	Schizophrenic		0.2043	0.7957		
22	Normal	Normal		0.8838	0.1162		
23	Schizophrenic	Schizophrenic		0.0591	0.9409		
24	Normal	Normal		0.9567	0.0433		
25	Schizophrenic	Normal	*	0.5513	0.4487		
26	Schizophrenic	Schizophrenic		0.2765	0.7235		
27	Normal	Normal		0.7258	0.2742		
28	Normal	Normal		0.9776	0.0224		
29	Normal	Normal		0.8097	0.1903		
30	Schizophrenic	Schizophrenic		0.0758	0.9242		
31	Schizophrenic	Schizophrenic		0.0120	0.9880		

* Misclassified

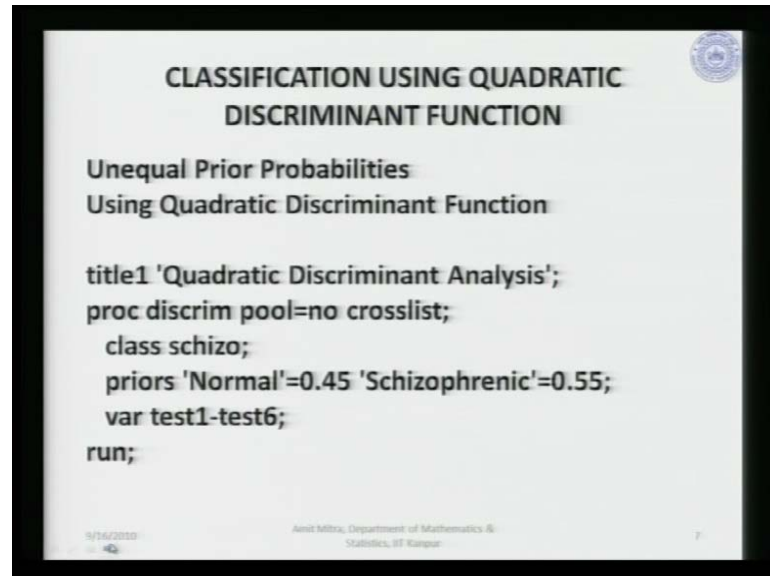
So, schizophrenic was a class which was having the class membership as 1. So, we have got these are the observations, these are the observation numbers. And the individual 1, the observation number 1 is coming from the class schizophrenic, it has got the posterior probabilities of the of in schizophrenic population as 0.84 and the posterior probability of a normal population, normal individual population is 0.15.

So, we see that this posterior probability is higher and hence what we have this observation classified correctly in to the schizophrenic class here. So, same as the interpretation for each records, each row of the records here say for example, if you look at the fifth record here, the fifth case is coming from a normal category of individuals. Now the posterior probabilities are coming out as 0.95 for the normal class and 0.04 for the schizophrenic class and since this is higher we classify correctly in to the normal class.

So, all these are correct classifications up to this particular point here. Now in the case number 7, there is a schizophrenic patient it is what the class membership is schizophrenic. And we are the posterior probability of the normal individual category as 0.99, and in the schizophrenic category the posterior probability 0.0098. And since, this is higher, we classify it as coming from a normal category of patients. However, that is a misclassification. So, this classification is wrongly done out here. So, these are the cases wherein we have got the misclassifications. So, all these stars here indicate that the

observations are going to misclassify based such posterior probability, the computed posterior probabilities.

(Refer Slide Time: 44:56)



CLASSIFICATION USING QUADRATIC DISCRIMINANT FUNCTION

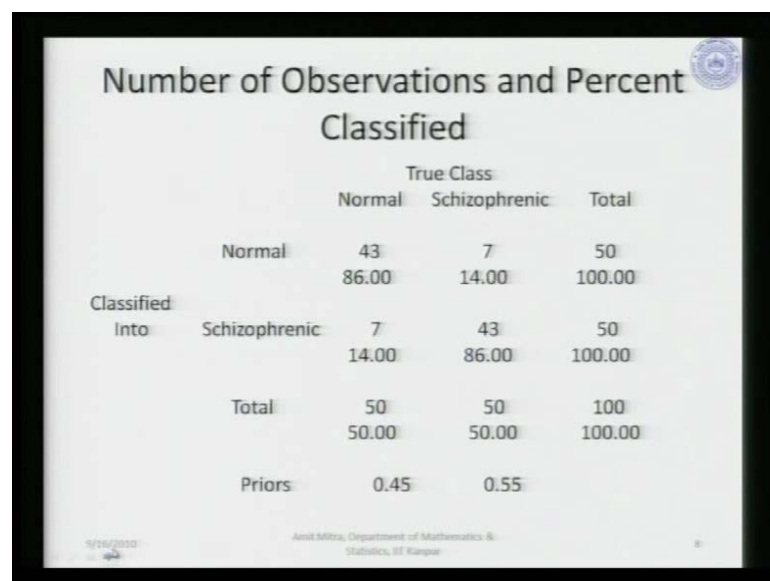
Unequal Prior Probabilities
Using Quadratic Discriminant Function

```
title1 'Quadratic Discriminant Analysis';  
proc discrim pool=no crosslist;  
  class schizo;  
  priors 'Normal'=0.45 'Schizophrenic'=0.55;  
  var test1-test6;  
run;
```

9/16/2010 Amit Mitra, Department of Mathematics & Statistics, IIT Kanpur 7

Next we look at the same data set only, and use a quadratic discriminant function with an unequal prior probabilities. So, this is what we are now looking at with unequal prior probabilities, and we are looking at a quadratic discriminant function. We take the 2 priors 0.45 and point 4 or 0.55.

(Refer Slide Time: 45:24)



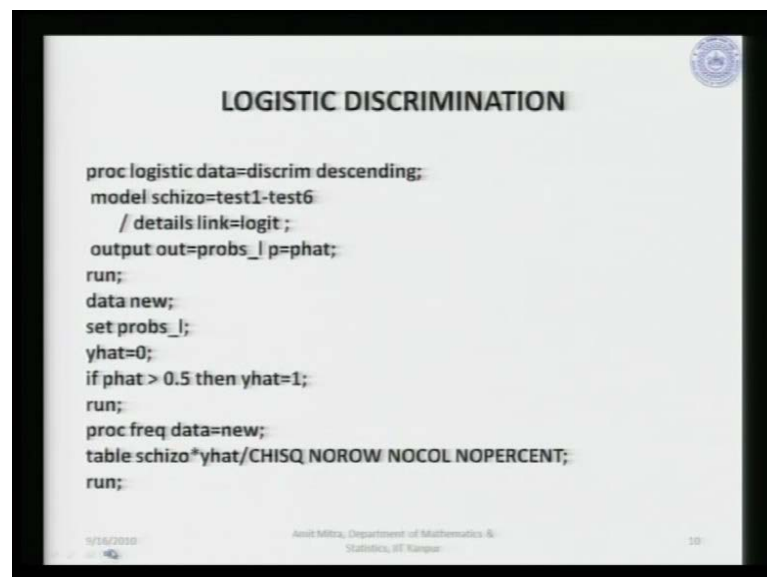
Number of Observations and Percent Classified

		True Class		
		Normal	Schizophrenic	Total
Classified Into	Normal	43 86.00	7 14.00	50 100.00
	Schizophrenic	7 14.00	43 86.00	50 100.00
Total		50 50.00	50 50.00	100 100.00
Priors		0.45	0.55	

9/16/2010 Amit Mitra, Department of Mathematics & Statistics, IIT Kanpur 8

So, these are the 2 prior probabilities. And with a quadratic discriminate function, we have the following confusion matrix along with the percent classifieds. So, this is once again, this is a true class membership which are 2 types normal, schizophrenic. And an individual coming from whichever class has got the possibility that it gets classified either in to normal or it is schizophrenic class.

Now, we see that from among 50 observations, now 43 in the normal category are correctly classified, 43 of them are also correctly classified. Then we have the number of observations coming here 50 from this class. And we have these numbers which are this classification numbers. So, these are the 2 quantities, wherein we have got these classifications. So, the percent of this classification here as we see its 14 percent out of this total number of cases. The similar is the interpretation, when we look at once again the posterior probability membership of membership in to the class schizophrenic.



The image shows a slide titled "LOGISTIC DISCRIMINATION" with SAS code. The code performs a logistic regression analysis on a dataset named 'discrim' using a descending method. It defines a model for 'schizo' with a logit link function. The output is saved to 'probs_1.p' as 'phat'. A new dataset 'new' is created with 'probs_1' and a predicted variable 'yhat' is set to 0, with a rule that if 'phat' is greater than 0.5, 'yhat' is set to 1. Finally, a frequency table is generated for 'schizo*yhat' with various statistics including CHISQ, NOROW, NOCOL, and NOPERCENT.

```

LOGISTIC DISCRIMINATION

proc logistic data=discrim descending;
  model schizo=test1-test6
    / details link=logit ;
  output out=probs_1 p=phat;
run;
data new;
  set probs_1;
  yhat=0;
  if phat > 0.5 then yhat=1;
run;
proc freq data=new;
  table schizo*yhat/CHISQ NOROW NOCOL NOPERCENT;
run;

```

5/16/2010
Anil Mitta, Department of Mathematics & Statistics, IIT Kanpur

Now, we have once again based on such posterior probabilities the classifications, we have misclassifications in some cases, total 17 in all. So, these are misclassification (()) the others are getting correctly classified using this quadratic discriminant function. We also apply a logistic discriminant function, and we look at what does the logistic discrimination that we learnt in today's lecture is going to lead us to. Now we use the proc logistic of the SAS procedures, in order to give this particular exercise with a link function as a logit link function, the type of link function that we have just now discussed, it is going to be that $\log \left(\frac{p}{1-p} \right)$ given x divided by $1 - \log \left(\frac{p}{1-p} \right)$ given x .

So, that is the with a logit link function, we are going to have these being classified. We take a cut off probability for a classification as 0.5. So, if the predictive probabilities are greater than 0.5, we take \hat{y} the predicted class membership to be equal to 1. And if it is otherwise, we take the predicted class membership to be equal to 0. This is associated with the schizophrenic class, and this is associated with the normal class. Now further more, once we have the classification done up to this particular point then the classification is done, we will look at the confusion matrix. And then, we will look at applying the a frequency procedure in order to look at the measures of association between such predicted memberships and the actual memberships.

(Refer Slide Time: 48:05)

Model Fit Statistics

Criterion	Intercept and Covariates	
	Only	
AIC	140.629	93.040
SC	143.235	111.277
-2 Log L	138.629	79.040

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	59.5890	6	<.0001
Score	46.5123	6	<.0001
Wald	26.2129	6	0.0002

9/16/2010 Amit Mitra, Department of Mathematics & Statistics, IIT Kanpur 11

So, these are elementary model fit statistic, this is the a L C S C Schwarz criterion minus two log l criterion. So, the log likelihood type of framework, now these are the hypothesis testing for the logistic regression setup, wherein we will we are testing beta equal to 0. That is the hypothesis of interest these are the various tests likelihood ratio test, the score test and the welds test. Each of them giving us a probability, which is very small less than 0.5.

(Refer Slide Time: 48:39)

The screenshot displays the output of a logistic regression analysis. The title is 'The LOGISTIC Procedure' and the subtitle is 'Analysis of Maximum Likelihood Estimates'. The main table shows the following data:

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	13.8733	2.9360	22.3286	<.0001
test1	1	-0.1598	0.0504	10.0398	0.0015
test2	1	-0.0461	0.0327	1.9880	0.1586
test3	1	-0.1138	0.0505	5.0756	0.0243
test4	1	0.0726	0.0454	2.5602	0.1096
test5	1	-0.0282	0.0404	0.4891	0.4843
test6	1	-0.00009	0.0391	0.0000	0.9981

Below this table, the 'Association of Predicted Probabilities and Observed Responses' is shown:

Percent Concordant	90.6	Somers' D	0.814
Percent Discordant	9.2	Gamma	0.815
Percent Tied	0.1	Tau-a	0.411
Pairs	2500	c	0.907

At the bottom of the slide, it says '9/16/2010' on the left, 'Asst. Prof., Department of Mathematics & Statistics, UT Kumpur' in the center, and '12' on the right.

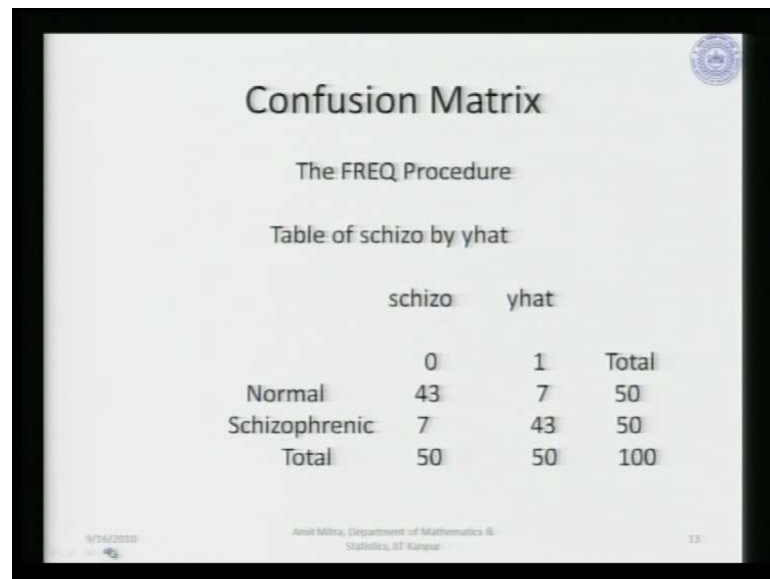
And what we have is these are the coefficients corresponding to such logistic regression of each of these tests, we assume that there is a constant term present.

So, these are the parameter, where the parameters are estimates **this** just shift a little bit, this is not standard estimate. This is to be taken with this so, it is standard error and welds chi square. So, this has to be coupled with this one. So, these are the parameter estimates, what we get from the data using an I R L S. And these are the corresponding weld the standard error quantities. So, these are the standard error columns and the welds chi square are these quantities. For those, which are less than a particular desired level of significance say 0.5, we reject the null hypothesis for all those terms there. So, the null hypothesis that this is equal to 0 is rejected, this is rejected this at 2 percent, so this is rejected. All other hypothesis are accepted at a 5 percent level of confidence.

Now, we also look at the association of the predicted probabilities through what when we are looking at the predicted class memberships as in after the predicted probabilities, if the predicted probability is greater than 0.5, we classify it in to say having the value equal to 1, y equal to 1 and 0 if it is otherwise. And then once again, we will be getting a confusion matrix. And from there, we look at the association of the predicted probabilities and the observed responses these are some standard measures of association, this is the percent concordant data. **Ah** after we have done the classification, this is a percent discordant in the data.

This is just an tied is 0.1 percent. So, we have a good fit, actually giving us percent concordant to be 90.6. So, these add up to 99.9, there is some round off somewhere. So, that 100 percent is not coming from all these cases. There are these many pairs, the Some are the criterion for measure of association is high, the gamma coefficient is high, tau a coefficient is high, the c kappa coefficient is also pretty high. So, the two the predicted probabilities and the observed responses, we naturally require them to be highly associated in order to have the classification to be worthwhile.

(Refer Slide Time: 51:00)



The slide displays a confusion matrix titled "Confusion Matrix" under the heading "The FREQ Procedure". The table is labeled "Table of schizo by yhat". The columns are labeled "schizo" and "yhat", and the rows are labeled "Normal", "Schizophrenic", and "Total". The matrix shows 43 correct classifications for the normal category and 43 correct classifications for the schizophrenic category. There are 7 misclassifications in each direction. The total number of observations is 100.

	schizo	yhat	Total
Normal	43	7	50
Schizophrenic	7	43	50
Total	50	50	100

9/16/2010 Amit Mitra, Department of Mathematics & Statistics, IIT Kanpur 13

And we have that here, now this is the confusion matrix what we have in here, it gives us once again 43 correct classifications from the normal category, and 43 correct classifications from the schizophrenic category. And these are the observations which are wrongly classified coming from the two different classes. This is the y hat quantity and this is the actual quantities, actual class memberships.

(Refer Slide Time: 51:27)

ACTUAL CLASS, PREDICTED PROBABILITY & PREDICTED CLASS			
Obs	class	phat	yhat
52	Schizophrenic	0.84029	1
53	Schizophrenic	0.56445	1
54	Schizophrenic	0.95652	1
55	Normal	0.03065	0
56	Normal	0.23743	0
57	Schizophrenic	0.71171	1
58	Schizophrenic	0.88218	1
59	Normal	0.55608	1
60	Normal	0.34181	0
61	Normal	0.00593	0
62	Schizophrenic	0.98746	1
63	Normal	0.02652	0
64	Normal	0.41416	0
65	Normal	0.53763	1
66	Normal	0.01952	0
67	Schizophrenic	0.16151	0
68	Schizophrenic	0.84072	1
69	Normal	0.05923	0
70	Schizophrenic	0.95532	1
71	Schizophrenic	0.43742	0
72	Schizophrenic	0.54458	1
73	Normal	0.64611	1
74	Normal	0.00368	0
75	Normal	0.04899	0
76	Schizophrenic	0.52340	1
77	Normal	0.25171	0
78	Schizophrenic	0.91828	1
79	Schizophrenic	0.94677	1
80	Normal	0.09721	0
81	Normal	0.05412	0
82	Schizophrenic	0.98622	1
83	Schizophrenic	0.96298	1
84	Schizophrenic	0.89956	1
85	Schizophrenic	0.95113	1
86	Schizophrenic	0.62963	1
87	Normal	0.01627	0
88	Normal	0.24129	0
89	Normal	0.22175	0
90	Schizophrenic	0.88076	1
91	Schizophrenic	0.57400	1
92	Schizophrenic	0.40999	0
93	Schizophrenic	0.95839	1

So, what we have here is just a representative of the data after the model has been fitted. So, this is an observation number 52, it is a schizophrenic class membership. We have the predicted probability. Now, what is this predicted probability, this is probability that y is equal to 1, that is it is schizophrenic given x . So, that probability is 0.8402. So, it is greater than 0.5 and hence the predicted class membership is given as 1.

Now, this is a correct classification. Similarly, we have all these predicted probabilities, which is probability of y equal to 1 y_i rather, corresponding to this case y_i equal to 1 given x quantities, and accordingly we have these. In situations, where this predicted probability is less than or equal to 0.5 as in this particular case, we will have that being classified in to the y hat category as taking the value 0.

So, that is in the normal category of patients. So, this is correctly classified, this is correctly classified. This also are correct classifications, correct classifications; however, this is a wrong classification, because we have this probability greater than 0.5. We classified that as y equal to 1 that is a predicted class membership is schizophrenic, which is wrong because the actual membership is normal. And hence, we have all other records in a similar way. Now for a same data set, we also apply a nearest neighbor classifier that we have learnt.

So, once again a proc discriminant, disc rim is used from SAS procedures with this non parametric method. Because it is a non parametric method, it is a nearest neighbor

classifier; we use the same data and get results. So, these are the nearest neighbor classifier examples with Euclidean distances, these are the cross validation results. It is interesting to look at what is the confusion matrix, and how it is behaving. So, we see that, this is what we have the true class membership 41 out of 49 are correctly classified, 44 out of 51 are correctly classified, and these are wrong classifications using a nearest neighbor classifier. Now, we have a multiclass problem, maybe we will take it in the next lecture.